

Proximal Quasi-Newton for Computationally Intensive ℓ_1 -regularized M -estimators

Kai Zhong¹ Ian E.H. Yen² Inderjit S. Dhillon² Pradeep Ravikumar²
¹ Institute for Computational Engineering & Sciences ² Department of Computer Science
 University of Texas at Austin
 zhongkai@ices.utexas.edu, {ianyeny, inderjit, pradeepr}@cs.utexas.edu

January 26, 2015

Abstract

We consider the class of optimization problems arising from computationally intensive ℓ_1 -regularized M -estimators, where the function or gradient values are very expensive to compute. A particular instance of interest is the ℓ_1 -regularized MLE for learning Conditional Random Fields (CRFs), which are a popular class of statistical models for varied structured prediction problems such as sequence labeling, alignment, and classification with label taxonomy. ℓ_1 -regularized MLEs for CRFs are particularly expensive to optimize since computing the gradient values requires an expensive inference step. In this work, we propose the use of a carefully constructed proximal quasi-Newton algorithm for such computationally intensive M -estimation problems, where we employ an aggressive active set selection technique. In a key contribution of the paper, we show that the proximal quasi-Newton method is provably *super-linearly convergent*, even in the absence of strong convexity, by leveraging a restricted variant of strong convexity. In our experiments, the proposed algorithm converges considerably faster than current state-of-the-art on the problems of sequence labeling and hierarchical classification.

1 Introduction

ℓ_1 -regularized M -estimators have attracted considerable interest in recent years due to their ability to fit large-scale statistical models, where the underlying model parameters are sparse. The optimization problem underlying these ℓ_1 -regularized M -estimators takes the form:

$$\min_{\mathbf{w}} f(\mathbf{w}) := \lambda \|\mathbf{w}\|_1 + \ell(\mathbf{w}), \quad (1)$$

where $\ell(\mathbf{w})$ is a convex differentiable loss function. In this paper, we are particularly interested in the case where the function or gradient values are very expensive to compute; we refer to these functions as computationally intensive functions, or **CI** functions in short. A particular case of interest are ℓ_1 -regularized MLEs for Conditional Random Fields (CRFs), where computing the gradient requires an expensive inference step.

There has been a line of recent work on computationally efficient methods for solving (1), including [2, 8, 13, 21, 23, 4]. It has now become well understood that it is key to leverage the sparsity of the optimal solution by maintaining sparse intermediate iterates [2, 5, 8]. Coordinate Descent (CD) based methods, like CDN [8], maintain the sparsity of intermediate iterates by focusing on an active set of working variables. A caveat with such methods is that, for CI functions, each coordinate update typically requires a call of inference oracle to evaluate partial derivative for single coordinate. One approach adopted in [16] to address this is using Blockwise Coordinate Descent that updates a block of variables at a time by ignoring the second-order effect, which however sacrifices the convergence guarantee. Newton-type methods have also attracted a surge of interest in recent years [5, 13], but these require computing the exact Hessian or Hessian-vector product, which is very expensive for CI functions. This then suggests the use of quasi-Newton methods, popular instances of which include OWL-QN [23], which is adapted from ℓ_2 -regularized L-BFGS, as well as Projected Quasi-Newton (PQN) [4]. A key caveat with OWL-QN and PQN however is that they do not exploit the sparsity of the underlying solution. In this paper, we consider the class of *Proximal Quasi-Newton* (Prox-QN) methods, which we argue seem particularly well-suited to such CI functions, for the following three reasons. Firstly, it requires gradient evaluations only once in each outer iteration. Secondly, it is a second-order method, which has asymptotic superlinear convergence. Thirdly, it can

employ some active-set strategy to reduce the time complexity from $O(d)$ to $O(nnz)$, where d is the number of parameters and nnz is the number of non-zero parameters.

While there has been some recent work on Prox-QN algorithms [2, 3], we carefully construct an implementation that is particularly suited to CI ℓ_1 -regularized M -estimators. We carefully maintain the sparsity of intermediate iterates, and at the same time reduce the gradient evaluation time. A key facet of our approach is our aggressive active set selection (which we also term a "shrinking strategy") to reduce the number of active variables under consideration at any iteration, and correspondingly the number of evaluations of partial gradients in each iteration. Our strategy is particularly aggressive in that it runs over multiple epochs, and in each epoch, chooses the next working set as a subset of the current working set rather than the whole set; while at the end of an epoch, allows for other variables to come in. As a result, in most iterations, our aggressive shrinking strategy only requires the evaluation of partial gradients in the current working set. Moreover, we adapt the L-BFGS update to the shrinking procedure such that the update can be conducted *without any loss of accuracy* caused by aggressive shrinking. Thirdly, we store our data in a *feature-indexed* structure to combine data sparsity as well as iterate sparsity.

[26] showed global convergence and asymptotic superlinear convergence for Prox-QN methods under the assumption that the loss function is *strongly convex*. However, this assumption is known to fail to hold in high-dimensional sampling settings, where the Hessian is typically rank-deficient, or indeed even in low-dimensional settings where there are redundant features. In a key contribution of the paper, we provide provable guarantees of asymptotic superlinear convergence for Prox-QN method, even without assuming strong-convexity, but under a restricted variant of strong convexity, termed Constant Nullspace Strong Convexity (CNSC), which is typically satisfied by standard M -estimators.

To summarize, our contributions are twofold. (a) We present a carefully constructed proximal quasi-Newton method for computationally intensive (CI) ℓ_1 -regularized M -estimators, which we empirically show to outperform many state-of-the-art methods on CRF problems. (b) We provide the first proof of asymptotic superlinear convergence for Prox-QN methods without strong convexity, but under a restricted variant of strong convexity, satisfied by typical M -estimators, including the ℓ_1 -regularized CRF MLEs.

2 Proximal Quasi-Newton Method

A proximal quasi-Newton approach to solve M -estimators of the form (1) proceeds by iteratively constructing a quadratic approximation of the objective function (1) to find the quasi-Newton direction, and then conducting a line search procedure to obtain the next iterate.

Given a solution estimate \mathbf{w}_t at iteration t , the proximal quasi-Newton method computes a descent direction by minimizing the following regularized quadratic model,

$$\mathbf{d}_t = \arg \min_{\Delta} \mathbf{g}_t^T \Delta + \frac{1}{2} \Delta^T B_t \Delta + \lambda \|\mathbf{w}_t + \Delta\|_1 \quad (2)$$

where $\mathbf{g}_t = \mathbf{g}(\mathbf{w}_t)$ is the gradient of $\ell(\mathbf{w}_t)$ and B_t is an approximation to the Hessian of $\ell(\mathbf{w})$. B_t is usually formulated by the L-BFGS algorithm. This subproblem (2) can be efficiently solved by randomized coordinate descent algorithm as shown in Section 2.2.

The next iterate is obtained from the backtracking line search procedure, $\mathbf{w}_{t+1} = \mathbf{w}_t + \alpha_t \mathbf{d}_t$, where the step size α_t is tried over $\{\beta^0, \beta^1, \beta^2, \dots\}$ until the Armijo rule is satisfied,

$$f(\mathbf{w}_t + \alpha_t \mathbf{d}_t) \leq f(\mathbf{w}_t) + \alpha_t \sigma \Delta_t,$$

where $0 < \beta < 1$, $0 < \sigma < 1$ and $\Delta_t = \mathbf{g}_t^T \mathbf{d}_t + \lambda(\|\mathbf{w}_t + \mathbf{d}_t\|_1 - \|\mathbf{w}_t\|_1)$.

2.1 BFGS update formula

B_t can be efficiently updated by the gradients of the previous iterations according to the BFGS update [18],

$$B_t = B_{t-1} - \frac{B_{t-1} \mathbf{s}_{t-1} \mathbf{s}_{t-1}^T B_{t-1}}{\mathbf{s}_{t-1}^T B_{t-1} \mathbf{s}_{t-1}} + \frac{\mathbf{y}_{t-1} \mathbf{y}_{t-1}^T}{\mathbf{y}_{t-1}^T \mathbf{s}_{t-1}} \quad (3)$$

where $\mathbf{s}_t = \mathbf{w}_{t+1} - \mathbf{w}_t$ and $\mathbf{y}_t = \mathbf{g}_{t+1} - \mathbf{g}_t$

We use the compact formula for B_t [18],

$$B_t = B_0 - Q R Q^T = B_0 - Q \hat{Q},$$

where

$$Q := \begin{bmatrix} B_0 S_t & Y_t \end{bmatrix}, R := \begin{bmatrix} S_t^T B_0 S_t & L_t \\ L_t^T & -D_t \end{bmatrix}^{-1}, \hat{Q} := R Q^T$$

$$S_t = [s_0, s_1, \dots, s_{t-1}], Y_t = [y_0, y_1, \dots, y_{t-1}]$$

$$D_t = \text{diag}[s_0^T y_0, \dots, s_{t-1}^T y_{t-1}] \text{ and } (L_t)_{i,j} = \begin{cases} s_{i-1}^T y_{j-1} & \text{if } i > j \\ 0 & \text{otherwise} \end{cases}$$

In practical implementation, we apply Limited-memory-BFGS. It only uses the information of the most recent m gradients, so that Q and \hat{Q} have only size, $d \times 2m$ and $2m \times d$, respectively. B_0 is usually set as $\gamma_t I$ for computing B_t , where $\gamma_t = y_{t-1}^T s_{t-1} / s_{t-1}^T s_{t-1}$ [18]. As will be discussed in Section 2.3, $Q(\hat{Q})$ is updated just on the rows(columns) corresponding to the working set, \mathcal{A} . The time complexity for L-BFGS update is $O(m^2|\mathcal{A}| + m^3)$.

2.2 Coordinate Descent for Inner Problem

Randomized coordinate descent is carefully employed to solve the inner problem (2) by Tang and Scheinberg [2]. In the update for coordinate j , $\mathbf{d} \leftarrow \mathbf{d} + z^* \mathbf{e}_j$, z^* is obtained by solving the one-dimensional problem,

$$z^* = \arg \min_z \frac{1}{2} (B_t)_{jj} z^2 + ((\mathbf{g}_t)_j + (B_t \mathbf{d})_j) z + \lambda |(\mathbf{w}_t)_j + d_j + z|$$

This one-dimensional problem has a closed-form solution, $z^* = -c + \mathcal{S}(c - b/a, \lambda/a)$, where \mathcal{S} is the soft-threshold function and $a = (B_t)_{jj}$, $b = (\mathbf{g}_t)_j + (B_t \mathbf{d})_j$ and $c = (\mathbf{w}_t)_j + d_j$. For $B_0 = \gamma_t I$, the diagonal of B_t can be computed by $(B_t)_{jj} = \gamma_t - \mathbf{q}_j^T \hat{\mathbf{q}}_j$, where \mathbf{q}_j^T is the j -th row of Q and $\hat{\mathbf{q}}_j$ is the j -th column of \hat{Q} . And the second term in b , $(B_t \mathbf{d})_j$ can be computed by,

$$(B_t \mathbf{d})_j = \gamma_t d_j - \mathbf{q}_j^T \hat{Q} \mathbf{d} = \gamma_t d_j - \mathbf{q}_j^T \hat{\mathbf{d}},$$

where $\hat{\mathbf{d}} := \hat{Q} \mathbf{d}$. Since $\hat{\mathbf{d}}$ has only $2m$ dimension, it is fast to update $(B_t \mathbf{d})_j$ by \mathbf{q}_j and $\hat{\mathbf{d}}$. In each inner iteration, only d_j is updated, so we have the fast update of $\hat{\mathbf{d}}$, $\hat{\mathbf{d}} \leftarrow \hat{\mathbf{d}} + \hat{\mathbf{q}}_j z^*$.

Since we only update the coordinates in the working set, the above algorithm has only computation complexity $O(m|\mathcal{A}| \times \text{inner_iter})$, where inner_iter is the number of iterations used for solving the inner problem.

2.3 Implementation

In this section, we discuss several key implementation details used in our algorithm to speed up the optimization. **Shrinking Strategy**

In each iteration, we select an active or working subset \mathcal{A} of the set of all variables: only the variables in this set are updated in the current iteration. The complementary set, also called the fixed set, has only values of zero and is not updated. The use of such a shrinking strategy reduces the overall complexity from $O(d)$ to $O(|\mathcal{A}|)$. Specifically, we (a) update the gradients just on the working set, (b) update $Q(\hat{Q})$ just on the rows(columns) corresponding to the working set, and (c) compute the latest entries in D_t, γ_t, L_t and $S_t^T S_t$ by just using the corresponding working set rather than the whole set.

The key facet of our “shrinking strategy” however is in aggressively shrinking the active set: at the next iteration, we set the active set to be a subset of the previous active set, so that $\mathcal{A}_t \subset \mathcal{A}_{t-1}$. Such an aggressive shrinking strategy however is not guaranteed to only weed out irrelevant variables. Accordingly, we proceed in epochs. In each epoch, we progressively shrink the active set as above, till the iterations seem to converge. At that time, we then allow for all the “shrunk” variables to come back and start a new epoch. Such a strategy was also called an ϵ -cooling strategy by Fan et al. [14], where the shrinking stopping criterion is loose at the beginning, and progressively becomes more strict each time all the variables are brought back. For L-BFGS update, when a new epoch starts, the memory of L-BFGS is cleaned to prevent any loss of accuracy.

Because at the first iteration of each new epoch, the entire gradient over all coordinates is evaluated, the computation time for those iterations accounts for a significant portion of the total time complexity. Fortunately, our experiments show that the number of epochs is typically between 3-5.

Inexact inner problem solution

Like many other proximal methods, e.g. GLMNET and QUIC, we solve the inner problem inexactly. This reduces the time complexity of the inner problem dramatically. The amount of inexactness is based on a heuristic method

which aims to balance the computation time of the inner problem in each outer iteration. The computation time of the inner problem is determined by the number of inner iterations and the size of working set. Thus, we let the number of inner iterations, $inner_iter = \min\{max_inner, \lfloor d/|\mathcal{A}| \rfloor\}$, where $max_inner = 10$ in our experiment.

Data Structure for both model sparsity and data sparsity

In our implementation we take two sparsity patterns into consideration: (a) model sparsity, which accounts for the fact that most parameters are equal to zero in the optimal solution; and (b) data sparsity, wherein most feature values of any particular instance are zeros. We use a *feature-indexed* data structure to take advantage of both sparsity patterns. Computations involving data will be time-consuming if we compute over all the instances including those that are zero. So we leverage the sparsity of data in our experiment by using vectors of pairs, whose members are the index and its value. Traditionally, each vector represents an instance and the indices in its pairs are the feature indices. However, in our implementation, to take both model sparsity and data sparsity into account, we use an inverted data structure, where each vector represents one feature (*feature-indexed*) and the indices in its pairs are the instance indices. This data structure facilitates the computation of the gradient for a particular feature, which involves only the instances related to this feature.

We summarize these steps in the algorithm below. And a detailed algorithm is in Appendix B.

Algorithm 1 Proximal Quasi-Newton Algorithm (Prox-QN)

Input: Dataset $\{\mathbf{x}^{(i)}, \mathbf{y}^{(i)}\}_{i=1,2,\dots,N}$, termination criterion ϵ , λ and L-BFGS memory size m .

Output: \mathbf{w}^* converging to $\arg \min_{\mathbf{w}} f(\mathbf{w})$.

- 1: Initialize $\mathbf{w} \leftarrow \mathbf{0}$, $\mathbf{g} \leftarrow \partial \ell(\mathbf{w})/\partial \mathbf{w}$, working set $\mathcal{A} \leftarrow \{1, 2, \dots, d\}$, and $S, Y, Q, \hat{Q} \leftarrow \phi$.
 - 2: **while** termination criterion is not satisfied or working set doesn't contain all the variables **do**
 - 3: Shrink working set.
 - 4: **if** Shrinking stopping criterion is satisfied **then**
 - 5: Take all the shrunken variables back to working set and clean the memory of L-BFGS.
 - 6: Update Shrinking stopping criterion and continue.
 - 7: **end if**
 - 8: Solve inner problem (2) over working set and obtain the new direction \mathbf{d} .
 - 9: Conduct line search based on Armijo rule and obtain new iterate \mathbf{w} .
 - 10: Update \mathbf{g} , S , Y , Q , \hat{Q} and related matrices over working set.
 - 11: **end while**
-

3 Convergence Analysis

In this section, we analyze the convergence behavior of proximal quasi-Newton method in the super-linear convergence phase, where the unit step size is chosen. To simplify the analysis, in this section, we assume the inner problem is solved exactly and no shrinking strategy is employed. We also provide the global convergence proof for Prox-QN method with shrinking strategy in Appendix A.5. In current literature, the analysis of proximal Newton-type methods relies on the assumption of *strongly convex* objective function to prove superlinear convergence [3]; otherwise, only sublinear rate can be proved [25]. However, our objective (1) is not strongly convex when the dimension is very large or there are redundant features. In particular, the Hessian matrix $H(\mathbf{w})$ of the smooth function $\ell(\mathbf{w})$ is not positive-definite. We thus leverage a recently introduced restricted variant of strong convexity, termed Constant Nullspace Strong Convexity (CNSC) in [1]. There the authors analyzed the behavior of proximal gradient and proximal Newton methods under such a condition. The proximal *quasi-Newton* procedure in this paper however requires a subtler analysis, but in a key contribution of the paper, we are nonetheless able to show asymptotic superlinear convergence of the Prox-QN method under this restricted variant of strong convexity.

Definition 1 (Constant Nullspace Strong Convexity (CNSC)). *A composite function (1) is said to have Constant Nullspace Strong Convexity restricted to space \mathcal{T} (CNSC- \mathcal{T}) if there is a constant vector space \mathcal{T} s.t. $\ell(\mathbf{w})$ depends only on $\text{proj}_{\mathcal{T}}(\mathbf{w})$, i.e. $\ell(\mathbf{w}) = \ell(\text{proj}_{\mathcal{T}}(\mathbf{w}))$, and its Hessian satisfies*

$$m\|\mathbf{v}\|^2 \leq \mathbf{v}^T H(\mathbf{w})\mathbf{v} \leq M\|\mathbf{v}\|^2, \quad \forall \mathbf{v} \in \mathcal{T}, \forall \mathbf{w} \in \mathbb{R}^d \quad (4)$$

for some $M \geq m > 0$, and

$$H(\mathbf{w})\mathbf{v} = \mathbf{0}, \quad \forall \mathbf{v} \in \mathcal{T}^\perp, \forall \mathbf{w} \in \mathbb{R}^d, \quad (5)$$

where $\text{proj}_{\mathcal{T}}(\mathbf{w})$ is the projection of \mathbf{w} onto \mathcal{T} and \mathcal{T}^\perp is the complementary space orthogonal to \mathcal{T} .

This condition can be seen to be an algebraic condition that is satisfied by typical M -estimators considered in high-dimensional settings. In this paper, we will abuse the use of CNSC- \mathcal{T} for symmetric matrices. We say a symmetric matrix H satisfies CNSC- \mathcal{T} condition if H satisfies (4) and (5). In the following theorems, we will denote the orthogonal basis of \mathcal{T} as $U \in \mathbb{R}^{d \times \hat{d}}$, where $\hat{d} \leq d$ is the dimensionality of \mathcal{T} space and $U^T U = I$. Then the projection to \mathcal{T} space can be written as $\text{proj}_{\mathcal{T}}(\mathbf{w}) = U U^T \mathbf{w}$.

Theorem 1 (Asymptotic Superlinear Convergence). *Assume $\nabla^2 \ell(\mathbf{w})$ and $\nabla \ell(\mathbf{w})$ are Lipschitz continuous. Let B_t be the matrices generated by BFGS update (3). Then if $\ell(\mathbf{w})$ and B_t satisfy CNSC- \mathcal{T} condition, the proximal quasi-Newton method has q -superlinear convergence:*

$$\|\mathbf{z}_{t+1} - \mathbf{z}^*\| \leq o(\|\mathbf{z}_t - \mathbf{z}^*\|),$$

where $\mathbf{z}_t = U^T \mathbf{w}_t$, $\mathbf{z}^* = U^T \mathbf{w}^*$ and \mathbf{w}^* is an optimal solution of (1).

The proof is given in Appendix A.4. We prove it by exploiting the CNSC- \mathcal{T} property. First, we re-build our problem and algorithm on the reduced space $\mathcal{Z} = \{\mathbf{z} \in \mathbb{R}^{\hat{d}} | \mathbf{z} = U^T \mathbf{w}\}$, where the strong-convexity property holds. Then we prove the asymptotic superlinear convergence on \mathcal{Z} following Theorem 3.7 in [26].

Theorem 2. *For Lipschitz continuous $\ell(\mathbf{w})$, the sequence $\{\mathbf{w}_t\}$ produced by the proximal quasi-Newton Method in the super-linear convergence phase has*

$$f(\mathbf{w}_t) - f(\mathbf{w}^*) \leq L \|\mathbf{z}_t - \mathbf{z}^*\|, \quad (6)$$

where $L = L_\ell + \lambda\sqrt{\hat{d}}$, L_ℓ is the Lipschitz constant of $\ell(\mathbf{w})$, $\mathbf{z}_t = U^T \mathbf{w}_t$ and $\mathbf{z}^* = U^T \mathbf{w}^*$.

The proof is also in Appendix A.4. It is proved by showing that both the smooth part and the non-differentiable part satisfy the modified Lipschitz continuity.

4 Application to Conditional Random Fields with ℓ_1 Penalty

In CRF problems, we are interested in learning a conditional distribution of labels $\mathbf{y} \in \mathcal{Y}$ given observation $\mathbf{x} \in \mathcal{X}$, where \mathbf{y} has application-dependent structure such as sequence, tree, or table in which label assignments have inter-dependency. The distribution is of the form

$$P_{\mathbf{w}}(\mathbf{y}|\mathbf{x}) = \frac{1}{Z_{\mathbf{w}}(\mathbf{x})} \exp \left\{ \sum_{k=1}^d w_k f_k(\mathbf{y}, \mathbf{x}) \right\},$$

where f_k is the feature functions, w_k is the associated weight, d is the number of feature functions and $Z_{\mathbf{w}}(\mathbf{x})$ is the partition function. Given a training data set $\{(\mathbf{x}_i, \mathbf{y}_i)\}_{i=1}^N$, our goal is to find the optimal weights \mathbf{w} such that the following ℓ_1 -regularized negative log-likelihood is minimized.

$$\min_{\mathbf{w}} f(\mathbf{w}) = \lambda \|\mathbf{w}\|_1 - \sum_{i=1}^N \log P_{\mathbf{w}}(\mathbf{y}^{(i)}|\mathbf{x}^{(i)}) \quad (7)$$

Since $|\mathcal{Y}|$, the number of possible values \mathbf{y} takes, can be exponentially large, the evaluation of $\ell(\mathbf{w})$ and the gradient $\nabla \ell(\mathbf{w})$ needs application-dependent oracles to conduct the summation over \mathcal{Y} . For example, in *sequence labeling problem*, a dynamic programming oracle, *forward-backward* algorithm, is usually employed to compute $\nabla \ell(\mathbf{w})$. Such an oracle can be very expensive. In Prox-QN algorithm for sequence labeling problem, the *forward-backward* algorithm takes $O(|\mathcal{Y}|^2 N T \times \text{exp})$ time, where *exp* is the time for the expensive exponential computation, T is the sequence length and \mathcal{Y} is the possible label set for a symbol in the sequence. Then given the obtained oracle, the evaluation of the partial gradients over the working set \mathcal{A} has time complexity, $O(D_{nnz}|\mathcal{A}|T)$, where D_{nnz} is the average number of instances related to a feature. Thus when $O(|\mathcal{Y}|^2 N T \times \text{exp} + D_{nnz}|\mathcal{A}|T) > O(m^3 + m^2|\mathcal{A}|)$, the gradients evaluation time will dominate.

The following theorem gives that the ℓ_1 -regularized CRF MLEs satisfy the CNSC- \mathcal{T} condition.

Theorem 3. *With ℓ_1 penalty, the CRF loss function, $\ell(\mathbf{w}) = -\sum_{i=1}^N \log P_{\mathbf{w}}(\mathbf{y}^{(i)}|\mathbf{x}^{(i)})$, satisfies the CNSC- \mathcal{T} condition with $\mathcal{T} = \mathcal{N}^\perp$, where $\mathcal{N} = \{\mathbf{v} \in \mathbb{R}^d | \Phi^T \mathbf{v} = 0\}$ is a constant subspace of \mathbb{R}^d and $\Phi \in \mathbb{R}^{d \times (N|\mathcal{Y}|)}$ is defined as below,*

$$\Phi_{kn} = f_k(\mathbf{y}_l, \mathbf{x}^{(i)}) - E[f_k(\mathbf{y}, \mathbf{x}^{(i)})]$$

where $n = (i-1)|\mathcal{Y}| + l$, $l = 1, 2, \dots, |\mathcal{Y}|$ and E is the expectation over the conditional probability $P_{\mathbf{w}}(\mathbf{y}|\mathbf{x}^{(i)})$.

According to the definition of CNSC- \mathcal{T} condition, the ℓ_1 -regularized CRF MLEs don't satisfy the classical strong-convexity condition when \mathcal{N} has non-zero members, which happens in the following two cases: (i) the exponential representation is not minimal [27], i.e. for any instance i there exist a non-zero vector \mathbf{a} and a constant b_i such that $\langle \mathbf{a}, \phi(\mathbf{y}, \mathbf{x}^{(i)}) \rangle = b_i$, where $\phi(\mathbf{y}, \mathbf{x}) = [f_1(\mathbf{y}, \mathbf{x}^{(i)}), f_2(\mathbf{y}, \mathbf{x}^{(i)}), \dots, f_d(\mathbf{y}, \mathbf{x}^{(i)})]^T$; (ii) $d > N|\mathcal{Y}|$, i.e., the number of feature functions is very large. The first case holds in many problems, like the sequence labeling and hierarchical classification discussed in Section 6, and the second case will hold in high-dimensional problems.

5 Related Methods

There have been several methods proposed for solving ℓ_1 -regularized M -estimators of the form in (7). In this section, we will discuss these in relation to our method.

Orthant-Wise Limited-memory Quasi-Newton (OWL-QN) introduced by Andrew and Gao [23] extends L-BFGS to ℓ_1 -regularized problems. In each iteration, OWL-QN computes a generalized gradient called *pseudo-gradient* to determine the orthant and the search direction, then does a line search and a projection of the new iterate back to the orthant. Due to its fast convergence, it is widely implemented by many software packages, such as CRF++, CRFsuite and Wapiti. But OWL-QN does not take advantage of the model sparsity in the optimization procedure, and moreover Yu et al. [22] have raised issues with its convergence proof.

Stochastic Gradient Descent (SGD) uses the gradient of a single sample as the search direction at each iteration. Thus, the computation for each iteration is very fast, which leads to fast convergence at the beginning. However, the convergence becomes slower than the second-order method when the iterate is close to the optimal solution. Recently, an ℓ_1 -regularized SGD algorithm proposed by Tsuruoka et al. [21] is claimed to have faster convergence than OWL-QN. It incorporates ℓ_1 -regularization by using a cumulative ℓ_1 penalty, which is close to the ℓ_1 penalty received by the parameter if it had been updated by the true gradient. Tsuruoka et al. do consider data sparsity, i.e. for each instance, only the parameters related to the current instance are updated. But they too do not take the model sparsity into account.

Coordinate Descent (CD) and **Blockwise Coordinate Descent** (BCD) are popular methods for ℓ_1 -regularized problem. In each coordinate descent iteration, it solves an one-dimensional quadratic approximation of the objective function, which has a closed-form solution. It requires the second partial derivative with respect to the coordinate. But as discussed by Sokolovska et al., the exact second derivative in CRF problem is intractable. So they instead use an approximation of the second derivative, which can be computed efficiently by the same inference oracle queried for the gradient evaluation. However, pure CD is very expensive because it requires to call the inference oracle for the instances related to the current coordinate in each coordinate update. BCD alleviates this problem by grouping the parameters with the same \mathbf{x} feature into a block. Then each block update only needs to call the inference oracle once for the instances related to the current \mathbf{x} feature. However, it cannot alleviate the large number of inference oracle calls unless the data is very sparse such that every instance appears only in very few blocks.

Proximal Newton method has proven successful on problems of ℓ_1 -regularized logistic regression [13] and Sparse Invariance Covariance Estimation [5], where the Hessian-vector product can be cheaply re-evaluated for each update of coordinate. However, the Hessian-vector product for CI function like CRF requires the query of the inference oracle no matter how many coordinates are updated at a time [17], which then makes the coordinate update on quadratic approximation as expensive as coordinate update in the original problem. Our proximal quasi-Newton method avoids such problem by replacing Hessian with a low-rank matrix from BFGS update.

6 Numerical Experiments

We compare our approach, Prox-QN, with four other methods, Proximal Gradient (Prox-GD), OWL-QN [23], SGD [21] and BCD [16]. For OWL-QN, we directly use the OWL-QN optimizer developed by Andrew et al.¹, where we set the memory size as $m = 10$, which is the same as that in Prox-QN. For SGD, we implement the algorithm proposed by Tsuruoka et al. [21], and use cumulative ℓ_1 penalty with learning rate $\eta_k = \eta_0 / (1 + k/N)$, where k is the SGD iteration and N is the number of samples. For BCD, we follow Sokolovska et al. [16] but with three modifications. First, we add a line search procedure in each block update since we found it is required for convergence. Secondly, we apply shrinking strategy as discussed in Section 2.3. Thirdly, when the second derivative for some coordinate is less than 10^{-10} , we set it to be 10^{-10} because otherwise the lack of ℓ_2 -regularization in our problem setting will lead to a very large new iterate.

¹<http://research.microsoft.com/en-us/downloads/b1eb1016-1738-4bd5-83a9-370c9d498a03/>

We evaluate the performance of Prox-QN method on two problems, sequence labeling and hierarchical classification. In particular, we plot the relative objective difference $(f(\mathbf{w}_t) - f(\mathbf{w}^*)) / f(\mathbf{w}^*)$ and the number of non-zero parameters (on a log scale) against time in seconds. More experiment results, for example, the testing accuracy and the performance for different λ 's, are in Appendix E. All the experiments are executed on 2.8GHz Intel Xeon E5-2680 v2 Ivy Bridge processor with 1/4TB memory and Linux OS.

6.1 Sequence Labeling

In sequence labeling problems, each instance $(\mathbf{x}, \mathbf{y}) = \{(\mathbf{x}_t, y_t)\}_{t=1,2,\dots,T}$ is a sequence of T pairs of observations and the corresponding labels. Here we consider the optical character recognition (OCR) problem, which aims to recognize the handwriting words. The dataset² was preprocessed by Taskar et al. [19] and was originally collected by Kassel [20], and contains 6877 words (instances). We randomly divide the dataset into two part: training part with 6216 words and testing part with 661 words. The character label set Y consists of 26 English letters and the observations are characters which are represented by images of 16 by 8 binary pixels as shown in Figure 1(a). We use degree 2 pixels as the raw features, which means all pixel pairs are considered. Therefore, the number of raw features is $J = 128 \times 127/2 + 128 + 1$, including a bias. For degree 2 features, $x_{tj} = 1$ only when both pixels are 1 and otherwise $x_{tj} = 0$, where x_{tj} is the j -th raw feature of \mathbf{x}_t . For the feature functions, we use unigram feature functions $\mathbf{1}(y_t = y, x_{tj} = 1)$ and bigram feature functions $\mathbf{1}(y_t = y, y_{t+1} = y')$ with their associated weights, $\Theta_{y,j}$ and $\Lambda_{y,y'}$, respectively. So $\mathbf{w} = \{\Theta, \Lambda\}$ for $\Theta \in \mathbb{R}^{|Y| \times J}$ and $\Lambda \in \mathbb{R}^{|Y| \times |Y|}$ and the total number of parameters, $d = |Y|^2 + |Y| \times J = 215,358$. Using the above feature functions, the potential function can be specified as, $\tilde{P}_{\mathbf{w}}(\mathbf{y}, \mathbf{x}) = \exp \left\{ \langle \Lambda, \sum_{t=1}^T (\mathbf{e}_{y_t} \mathbf{x}_t^T) \rangle + \langle \Theta, \sum_{t=1}^{T-1} (\mathbf{e}_{y_t} \mathbf{e}_{y_{t+1}}^T) \rangle \right\}$, where $\langle \cdot, \cdot \rangle$ is the sum of element-wise product and $\mathbf{e}_y \in \mathbb{R}^{|Y|}$ is a unit vector with 1 at y -th entry and 0 at other entries. The gradient and the inference oracle are given in Appendix D.1.

In our experiment, λ is set as 100, which leads to a relative high testing accuracy and an optimal solution with a relative small number of non-zero parameters (see Appendix E.2). The learning rate η_0 for SGD is tuned to be 2×10^{-4} for best performance. In BCD, the unigram parameters are grouped into J blocks according to the \mathbf{x} features while the bigram parameters are grouped into one block. Our proximal quasi-Newton method can be seen to be much faster than the other methods.

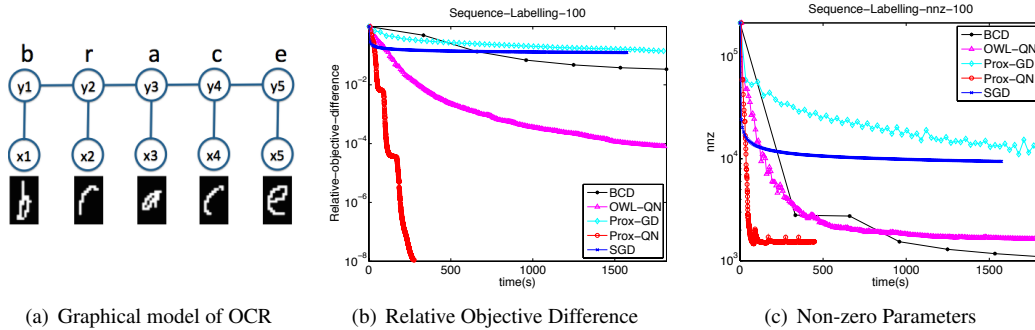


Figure 1: Sequence Labeling Problem

6.2 Hierarchical Classification

In hierarchical classification problems, we have a label taxonomy, where the classes are grouped into a tree as shown in Figure 2(a). Here $y \in \mathcal{Y}$ is one of the leaf nodes. If we have totally K classes (number of nodes) and J raw features, then the number of parameters is $d = K \times J$. Let $W \in \mathbb{R}^{K \times J}$ denote the weights. The feature function corresponding to $W_{k,j}$ is $f_{k,j}(y, \mathbf{x}) = \mathbf{1}[k \in \text{Path}(y)]x_j$, where $k \in \text{Path}(y)$ means class k is an ancestor of y or y itself. The potential function is $\tilde{P}_W(y, \mathbf{x}) = \exp \left\{ \sum_{k \in \text{Path}(y)} \mathbf{w}_k^T \mathbf{x} \right\}$ where \mathbf{w}_k^T is the weight vector of k -th class, i.e. the k -th row of W . The gradient and the inference oracle are given in Appendix D.2.

The dataset comes from Task1 of the dry-run dataset of LSHTC³. It has 4,463 samples, each with $J=51,033$ raw features. The hierarchical tree has 2,388 classes which includes 1,139 leaf labels. Thus, the number of the

²<http://www.seas.upenn.edu/~taskar/ocr/>

³<http://lshtc.iit.demokritos.gr/node/1>

parameters $d = 121,866,804$. The feature values are scaled by svm-scale program in the LIBSVM package. We set $\lambda = 1$ to achieve a relative high testing accuracy and high sparsity of the optimal solution. The SGD initial learning rate is tuned to be $\eta_0 = 10$ for best performance. In BCD, parameters are grouped into J blocks according to the raw features.

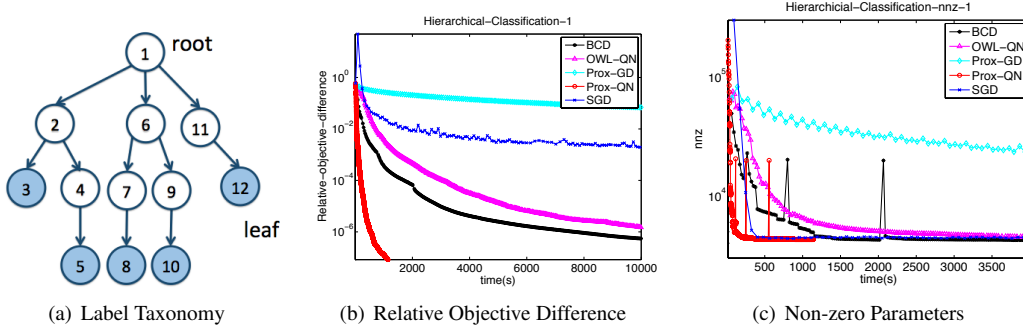


Figure 2: Hierarchical Classification Problem

As both Figure 1(b),1(c) and Figure 2(b),2(c) show, Prox-QN achieves much faster convergence and moreover obtains a sparse model in much less time.

Acknowledgement

This research was supported by NSF grants CCF-1320746 and CCF-1117055. P.R. acknowledges the support of ARO via W911NF-12-1-0390 and NSF via IIS-1149803, IIS-1320894, IIS-1447574, and DMS-1264033. K.Z. acknowledges the support of the National Initiative for Modeling and Simulation fellowship

References

- [1] I. E.H. Yen, C.-J. Hsieh, P. Ravikumar, and I. S. Dhillon. Constant Nullspace Strong Convexity and Fast Convergence of Proximal Methods under High-Dimensional Settings. In NIPS 2014.
- [2] X. Tang and K. Scheinberg. Efficiently Using Second Order Information in Large ℓ_1 Regularization Problems. arXiv:1303.6935, 2013.
- [3] J. D. Lee, Y. Sun, and M. A. Saunders. Proximal newton-type methods for minimizing composite functions. In NIPS 2012.
- [4] M. Schmidt, E. Van Den Berg, M.P. Friedlander, and K. Murphy. Optimizing costly functions with simple constraints: A limited-memory projected Quasi-Newton algorithm. In Int. Conf. Artif. Intell. Stat., 2009.
- [5] C.-J. Hsieh, M. A. Sustik, I. S. Dhillon, and P. Ravikumar. Sparse inverse covariance estimation using quadratic approximation. In NIPS 2011.
- [6] S. Boyd and L. Vandenberghe. Convex Optimization, Cambridge Univ. Press, Cambridge, U.K., 2003.
- [7] P.-W. Wang and C.-J. Lin. Iteration Complexity of Feasible Descent Methods for Convex Optimization. Technical report, Department of Computer Science, National Taiwan University, Taipei, Taiwan, 2013.
- [8] G.-X. Yuan, K.-W. Chang, C.-J. Hsieh, and C.-J. Lin. A comparison of optimization methods and software for large-scale ℓ_1 -regularized linear classification. Journal of Machine Learning Research (JMLR), 11:3183-3234, 2010.
- [9] A. Agarwal, S. Negahban, and M. Wainwright. Fast Global Convergence Rates of Gradient Methods for High-Dimensional Statistical Recovery. In NIPS 2010.
- [10] K. Hou, Z. Zhou, A. M.-S. So, and Z.-Q. Luo. On the linear convergence of the proximal gradient method for trace norm regularization. In NIPS 2014.

- [11] L. Xiao and T. Zhang. A proximal-gradient homotopy method for the ℓ_1 -regularized least-squares problem. In ICML 2012.
- [12] P. Tseng and S. Yun. A coordinate gradient descent method for nonsmooth separable minimization, *Math. Prog. B.* 117, 2009.
- [13] G.-X. Yuan, C.-H. Ho, and C.-J. Lin. An improved GLMNET for ℓ_1 -regularized logistic regression, *JMLR*, 13:1999-2030, 2012.
- [14] R.-E. Fan, K.-W. Chang, C.-J. Hsieh, X.-R. Wang, and C.-J. Lin. LIBLINEAR: A library for large linear classification, *JMLR*, 9:1871-1874, 2008.
- [15] A. J Hoffman. On approximate solutions of systems of linear inequalities. *Journal of Research of the National Bureau of Standards*, 1952.
- [16] N. Sokolovska, T. Lavergne, O. Cappe, and F. Yvon. Efficient Learning of Sparse Conditional Random Fields for Supervised Sequence Labelling. *arXiv:0909.1308*, 2009.
- [17] Y. Tsuboi, Y. Unno, H. Kashima, and N. Okazaki. Fast Newton-CG Method for Batch Learning of Conditional Random Fields, *Proceedings of the Twenty-Fifth AAAI Conference on Artificial Intelligence*, 2011.
- [18] J. Nocedal and S. J. Wright. *Numerical Optimization*. Springer Series in Operations Research. Springer, New York, NY, USA, 2nd edition, 2006.
- [19] B. Taskar, C. Guestrin, and D. Koller. Max-margin markov networks. In *NIPS 2003*.
- [20] R. Kassel. A Comparison of Approaches to On-line Handwritten Character Recognition. PhD thesis, MIT Spoken Language Systems Group, 1995.
- [21] Y. Tsuruoka, J. Tsujii, and S. Ananiadou. Stochastic gradient descent training for ℓ_1 -regularized log-linear models with cumulative penalty. In *Proceedings of the Joint Conference of the 47th Annual Meeting of the ACL and the 4th International Joint Conference on Natural Language Processing of the AFNLP*, pages 477-485, Suntec, Singapore, 2009.
- [22] J. Yu, S.V.N. Vishwanathan, S. Gunter, and N. N. Schraudolph. A Quasi-Newton approach to nonsmooth convex optimization problems in machine learning, *JMLR*, 11:1-57, 2010.
- [23] G. Andrew and J. Gao. Scalable training of ℓ_1 -regularized log-linear models. In *ICML 2007*.
- [24] J.E. Dennis and J.J. More. A characterization of superlinear convergence and its application to Quasi-Newton methods. *Math. Comp.*, 28(126):549560, 1974.
- [25] K. Scheinberg and X. Tang. Practical Inexact Proximal Quasi-Newton Method with Global Complexity Analysis. *COR@L Technical Report at Lehigh University*. *arXiv:1311.6547*, 2013
- [26] J. D. Lee, Y. Sun, and M. A. Saunders. Proximal Newton-type methods for minimizing composite functions. *arXiv:1206.1623*, 2012
- [27] M. J. Wainwright and M. I. Jordan. *Graphical models, exponential families, and variational inference*. Technical Report 649, Dept. Statistics, Univ. California, Berkeley. 2003

APPENDIX

A Convergence Proof

To exploit the CNSC- \mathcal{T} property, we first re-build our problem and algorithm on the reduced space $\mathcal{Z} = \{z \in \mathbb{R}^{\hat{d}} | z = U^T w\}$, where the strong-convexity property holds. Then we prove the asymptotic super-linear convergence on \mathcal{Z} under the condition that the inner problem is solved exactly and no shrinking strategy is not applied. Finally we prove the objective (1) is bounded by the difference between current iterate and the optimal solution. In Section A.5, we provide the global convergence proof when the shrinking strategy is applied.

A.1 Representing the problem in a reduced and compact space

Properties of CNSC- \mathcal{T} condition

For $\ell(w)$ satisfying CNSC- \mathcal{T} condition, we have $\ell(w) = \ell(\text{proj}_{\mathcal{T}}(w))$. Define g to be the gradient of $\ell(w)$ and H to be the Hessian of $\ell(w)$. As both g and H are in the \mathcal{T} space, we have $g(w) = UU^T g(\text{proj}_{\mathcal{T}}(w)) = g(\text{proj}_{\mathcal{T}}(w))$ and $H(w) = UU^T H(\text{proj}_{\mathcal{T}}(w))UU^T = H(\text{proj}_{\mathcal{T}}(w))$.

Objective formulation in the reduced space

Define $\hat{\ell}(z) = \ell(Uz)$. Then if $z = U^T w$, we have $\hat{\ell}(z) = \ell(w)$, $\hat{g}(z) = U^T g(w)$ and $\hat{H}(z) = U^T H(w)U$, where $\hat{g}(z)$ and $\hat{H}(z)$ are the gradient and Hessian of $\hat{\ell}(z)$ respectively. Now \hat{H} is positive definite with minimal eigenvalue m . The objective (1) can be re-formulated in the reduced space by

$$\min_z \hat{f}(z) = h(z) + \hat{\ell}(z), \quad (8)$$

where

$$h(z) = \min_{U^T w = z} \lambda \|w\|_1$$

We now prove that $h(z)$ is a convex function, i.e.,

$$ch(z_1) + (1-c)h(z_2) \geq h(cz_1 + (1-c)z_2)$$

for any $0 \leq c \leq 1$, z_1 and z_2 .

Proof. Let

$$w_1 = \underset{U^T w = z_1}{\operatorname{argmin}} \lambda \|w\|_1 \text{ and } w_2 = \underset{U^T w = z_2}{\operatorname{argmin}} \lambda \|w\|_1$$

Then,

$$\begin{aligned} ch(z_1) + (1-c)h(z_2) &= \lambda(c\|w_1\|_1 + (1-c)\|w_2\|_1) \\ &\geq \lambda(\|cw_1 + (1-c)w_2\|_1) \\ &\geq h(U^T(cw_1 + (1-c)w_2)) \\ &= h(cz_1 + (1-c)z_2) \end{aligned}$$

□

The optimal solution z^* of (8) has the following relationship with the optimal solution w^* of (1),

$$w^* = \underset{U^T w = z^*}{\operatorname{argmin}} \lambda \|w\|_1 \text{ and } z^* = U^T w^* \quad (9)$$

Lipschitz continuity in the reduced space

Throughout the paper, we assume the Hessian of $\ell(w)$ has Lipschitz continuity with constant L_H . According to the Lipschitz continuity,

$$\|H(w_2)(w_1 - w_2) - (g(w_1) - g(w_2))\| \leq \frac{L_H}{2} \|w_1 - w_2\|^2$$

In the corresponding reduced space, the Lipschitz continuity also holds with the same constant .

$$\|\hat{H}(z_2)(z_1 - z_2) - (\hat{g}(z_1) - \hat{g}(z_2))\| \leq \frac{L_H}{2} \|z_1 - z_2\|^2 \quad (10)$$

BFGS update formula in the reduced space

If B_0 is in the \mathcal{T} space, B_t is also in the \mathcal{T} space. This can be shown by re-formulating the BFGS update and mathematical induction,

$$B_t = U\hat{B}_{t-1}U^T - \frac{U\hat{B}_{t-1}U^T s_{t-1} s_{t-1}^T U\hat{B}_{t-1}U^T}{s_{t-1}^T U\hat{B}_{t-1}U^T s_{t-1}} + \frac{UU^T y_{t-1} y_{t-1}^T UU^T}{y_{t-1}^T UU^T s_{t-1}} \quad (11)$$

Thus

$$\hat{B}_t = \hat{B}_{t-1} - \frac{\hat{B}_{t-1} \hat{s}_{t-1} \hat{s}_{t-1}^T \hat{B}_{t-1}}{\hat{s}_{t-1}^T \hat{B}_{t-1} \hat{s}_{t-1}} + \frac{\hat{y}_{t-1} \hat{y}_{t-1}^T}{\hat{y}_{t-1}^T \hat{s}_{t-1}} \quad (12)$$

where $\hat{s} = U^T s$, $\hat{y} = U^T y$ and $U\hat{B}_t U^T = B_t$. It can be proved that \hat{B}_t generated in (12) is positive definite provided $\hat{y}^T \hat{s} > 0$ [18]. If we additionally assume $m\|z\|^2 \leq z^T \hat{B}_t z \leq M\|z\|^2$ for any $z \in \mathbb{R}^{\hat{d}}$, then B_t satisfies the CNSC- \mathcal{T} condition.

Iterate in the reduced space

The potential new iterate w^+ is

$$w^+ = \underset{v}{\operatorname{argmin}} \lambda \|v\|_1 + \frac{1}{2} (v - w_t)^T B_t (v - w_t) + g_t^T (v - w_t) \quad (13)$$

In the reduced space, the potential new iterate (13) can be represented by,

$$z^+ = \underset{x}{\operatorname{argmin}} h(x) + \frac{1}{2} (x - z_t)^T \hat{B}_t (x - z_t) + \hat{g}_t^T (x - z_t) \quad (14)$$

z^+ and w^+ also satisfy Equation (9), i.e.

$$w^+ = \underset{U^T w = z^+}{\operatorname{argmin}} \|w\|_1 \quad (15)$$

In this paper, we consider the convergence phase when z_t is close enough to the optimum such that the unit step size is always chosen, i.e. $z_{t+1} = z^+$ [26].

A.2 Global linear Convergence

Lemma 1 (Global linear Convergence). *For $\nabla \hat{\ell}(z)$ satisfying Lipschitz-continuity with a constant L_g and B_t satisfying CNSC- \mathcal{T} , the sequence $\{z_t\}_{t=1}^\infty$ produced by Prox-QN method converges at least R-linearly.*

Proof. This theorem follows Theorem 2 in [12], where the coordinate block J_k is chosen to be the whole coordinate set. Assumption 2(a) in [12] is satisfied because of Theorem 4 C4 in [12] by assuming $\nabla \hat{\ell}(z)$ is Lipschitz-continuous. Other conditions of Theorem 2 in [12] can be easily justified. \square

A.3 Quadratic Convergence of Proximal Newton Method and Dennis-More Criterion

Lemma 2 (Quadratic Convergence of Prox-Newton (Theorem 3 in [1])). *For $\ell(w)$ satisfying CNSC- \mathcal{T} with Lipschitz-continuous second derivative $H(w) = \nabla^2 \ell(w)$, the sequence $\{w_t\}$ produced by proximal Newton Method in the quadratic convergence phase has*

$$\|z_{t+1} - z^*\| \leq \frac{L_H}{2m} \|z_t - z^*\|^2,$$

where $z^* = U^T w^*$, $z_t = U^T w_t$, w^* is the optimal solution and L_H is the Lipschitz constant for $H(w)$.

Lemma 3. *If $B_0 = U\hat{B}_0 U^T$ satisfies CNSC- \mathcal{T} condition, then \hat{B}_t generated by (12) satisfies the Dennis-More criterion [24], namely,*

$$\lim_{t \rightarrow \infty} \frac{\|(\hat{B}_t - \hat{H}^*)(z_{t+1} - z_t)\|}{\|z_{t+1} - z_t\|} = 0,$$

where $\hat{H}^* = \nabla^2 \hat{\ell}(z^*)$ and z^* is the optimal solution of (8).

Proof. We want to show that this proof can follow the proof of Theorem 6.6 in [18]. We will verify that the conditions of Theorem 6.6 in [18] are satisfied here. First, the Lipschitz continuity of $\hat{H}(\mathbf{z})$ is implied by Lipschitz continuity of $H(\mathbf{w})$:

$$\begin{aligned}\|(\hat{H}(\mathbf{z}_1) - \hat{H}(\mathbf{z}_2))\| &= \|U^T(H(\mathbf{w}_1) - H(\mathbf{w}_2))U\| \\ &\leq \|H(\mathbf{w}_1) - H(\mathbf{w}_2)\| \\ &= \|H(U\mathbf{z}_1) - H(U\mathbf{z}_2)\| \\ &\leq L_H\|\mathbf{z}_1 - \mathbf{z}_2\|\end{aligned}$$

where the last inequality is from the Lipschitz continuity of $H(\mathbf{w})$. The second condition, $\sum_{t=0}^{\infty} \|\mathbf{z}_t - \mathbf{z}^*\| < \infty$, is implied by the global linear convergence(Lemma 1). \square

A.4 Asymptotic Superlinear Convergence

Proof of Theorem 1

Proof. If B_t satisfies CNSC- \mathcal{T} condition, then \hat{B}_t satisfies $m\|\mathbf{z}\|^2 \leq \mathbf{z}^T \hat{B}_t \mathbf{z} \leq M\|\mathbf{z}\|^2$ for any $\mathbf{z} \in \mathbb{R}^{\hat{d}}$. The Lipschitz-continuous H implies Lipschitz-continuity of \hat{H} . Therefore by applying the Prox-QN method in the reduced space, this theorem follows Theorem 3.7 in [26], Lemma 3 and Lemma 2. \square

Proof of Theorem 2

Proof. We prove this theorem by showing $|\ell(\mathbf{w}_t) - \ell(\mathbf{w}^*)| \leq L_\ell\|\mathbf{z}_t - \mathbf{z}^*\|$ and $\|\mathbf{w}_t\|_1 - \|\mathbf{w}^*\|_1 \leq \sqrt{d}\|\mathbf{z}_t - \mathbf{z}^*\|$. The first part is given by,

$$|\ell(\mathbf{w}_t) - \ell(\mathbf{w}^*)| = |\ell(UU^T\mathbf{w}_t) - \ell(UU^T\mathbf{w}^*)| \leq L_\ell\|UU^T(\mathbf{w}_t - \mathbf{w}^*)\| = L_\ell\|\mathbf{z}_t - \mathbf{z}^*\|$$

where the inequality comes from the Lipschitz-continuity of $\ell(\mathbf{w})$. In the super-linear convergence phase, the unit step size is chosen, so each iterate satisfies (15). We have $\|\mathbf{w}_t\|_1 \leq \|UU^T\mathbf{w}_t + (I - UU^T)\mathbf{w}^*\|_1$. Moreover, due to the Lipschitz-continuity of ℓ_1 norm, which is $\|\mathbf{w}\|_1 - \|\mathbf{v}\|_1 \leq \sqrt{d}\|\mathbf{w} - \mathbf{v}\|$, we have,

$$\begin{aligned}\|UU^T\mathbf{w}_t + (I - UU^T)\mathbf{w}^*\|_1 &\leq \|\mathbf{w}^*\|_1 + \sqrt{d}\|UU^T\mathbf{w}_t - UU^T\mathbf{w}^*\| \\ &\leq \|\mathbf{w}^*\|_1 + \sqrt{d}\|\mathbf{z}_t - \mathbf{z}^*\|\end{aligned}$$

\square

A.5 Global Convergence with Shrinking

In Theorem 1, we assume shrinking strategy is not employed and the inner problem is solved exactly. In this subsection, we show that by only assuming the inner problem is solved exactly, Prox-QN method with shrinking will still globally converge to the optimum under the CNSC- \mathcal{T} condition. We first prove that with sufficient small step size, the Armijo rule will be satisfied.

Lemma 4. *If the step size,*

$$\alpha \leq \min\left\{1, \frac{m}{L_1}(1 - \sigma)\right\}$$

then the Armijo rule is satisfied, i.e.,

$$f(\mathbf{w} + \alpha\mathbf{d}) \leq f(\mathbf{w}) + \alpha\sigma(\lambda\|\mathbf{w} + \mathbf{d}\|_1 - \lambda\|\mathbf{w}\|_1 + \mathbf{g}^T\mathbf{d})$$

where L_1 is the Lipschitz-continuity constant.

Proof. Let $\mathbf{w}^+ = \mathbf{w} + \alpha\mathbf{d}$,

$$\begin{aligned}f(\mathbf{w}^+) - f(\mathbf{w}) &= \ell(\mathbf{w}^+) - \ell(\mathbf{w}) + \lambda(\|\mathbf{w}^+\|_1 - \|\mathbf{w}\|_1) \\ &\leq \int_0^1 \nabla\ell(\mathbf{w} + s\alpha\mathbf{d})(\alpha\mathbf{d})ds + \alpha\lambda\|\mathbf{w} + \mathbf{d}\|_1 + (1 - \alpha)\lambda\|\mathbf{w}\|_1 - \lambda\|\mathbf{w}\|_1 \\ &= \alpha(\nabla\ell(\mathbf{w})^T\mathbf{d} + \lambda\|\mathbf{w} + \mathbf{d}\|_1 - \lambda\|\mathbf{w}\|_1) + \alpha \int_0^1 \mathbf{d}^T(\nabla\ell(\mathbf{w} + s\alpha\mathbf{d}) - \nabla\ell(\mathbf{w}))ds \\ &\leq \alpha(\nabla\ell(\mathbf{w})^T\mathbf{d} + \lambda\|\mathbf{w} + \mathbf{d}\|_1 - \lambda\|\mathbf{w}\|_1) + \alpha \int_0^1 \|U^T\mathbf{d}\|\|\nabla\ell(\mathbf{w} + s\alpha\mathbf{d}) - \nabla\ell(\mathbf{w})\|ds\end{aligned}$$

Because

$$\|\nabla\ell(\mathbf{w} + s\alpha\mathbf{d}) - \nabla\ell(\mathbf{w})\| = \|\nabla\ell(UU^T\mathbf{w} + s\alpha UU^T\mathbf{d}) - \nabla\ell(UU^T\mathbf{w})\| \leq sL_1\|U^T\mathbf{d}\|$$

we have

$$f(\mathbf{w}^+) - f(\mathbf{w}) \leq \alpha \left((\nabla\ell(\mathbf{w})^T\mathbf{d} + \lambda\|\mathbf{w} + \mathbf{d}\|_1 - \lambda\|\mathbf{w}\|_1) + \frac{L_1\alpha}{2}\|U^T\mathbf{d}\|^2 \right)$$

For $\alpha \leq \min\{1, \frac{m}{L_1}(1 - \sigma)\}$,

$$\frac{L_1\alpha}{2}\|U^T\mathbf{d}\|^2 \leq \frac{m}{2}(1 - \sigma)\|U^T\mathbf{d}\|^2 \leq \frac{1 - \sigma}{2}\mathbf{d}^T B \mathbf{d}$$

As \mathbf{d} minimizes Eq. (2) in the main paper, we have $\frac{1}{2}\mathbf{d}^T B \mathbf{d} \leq -(\nabla\ell(\mathbf{w})^T\mathbf{d} + \lambda\|\mathbf{w} + \mathbf{d}\|_1 - \lambda\|\mathbf{w}\|_1)$. So we obtain the sufficient descent condition,

$$f(\mathbf{w}^+) - f(\mathbf{w}) \leq \alpha\sigma (\nabla\ell(\mathbf{w})^T\mathbf{d} + \lambda\|\mathbf{w} + \mathbf{d}\|_1 - \lambda\|\mathbf{w}\|_1)$$

□

Proposition 1. Assume $\nabla^2\ell(\mathbf{w})$ and $\nabla\ell(\mathbf{w})$ are Lipschitz continuous. Let $\{B_t\}_{t=1,2,3,\dots}$ be the matrices generated by BFGS update. Then if $\ell(\mathbf{w})$ and B_t satisfy CNSC- \mathcal{T} condition and the inner problem is solved exactly, the proximal quasi-Newton method with shrinking has global convergence.

Proof. Our algorithm allows all the variables to re-enter the working set at the beginning of each epoch. And before it terminates all the variables must be checked. Thus as many as epochs are taken in the optimization procedure until the global stopping criterion is attained. Let's denote $\{t_k\}_{k=0,1,2,3,\dots}$ to be the iterations when an epochs begins. In these iterations, all the variables are taken into consideration. As shown in Lemma 4, there exists some constant α_0 ,

$$f(\mathbf{w}_{t_{k+1}}) - f(\mathbf{w}_{t_k}) \leq \alpha_0\sigma (\nabla\ell(\mathbf{w}_{t_k})^T\mathbf{d}_{t_k} + \lambda\|\mathbf{w}_{t_k} + \mathbf{d}_{t_k}\|_1 - \lambda\|\mathbf{w}_{t_k}\|_1)$$

And as in each epoch the function value is non-increasing across the iterations, i.e. for any k , $f(\mathbf{w}_{t_{k+1}}) \leq f(\mathbf{w}_{t_k+1})$. Thus, we have

$$f(\mathbf{w}_{t_{K+1}}) - f(\mathbf{w}_{t_0}) \leq \sum_{k=0}^K f(\mathbf{w}_{t_{k+1}}) - f(\mathbf{w}_{t_k}) \leq -\alpha_0\sigma \sum_{k=0}^K \mathbf{d}_{t_k}^T B_{t_k} \mathbf{d}_{t_k}$$

As $f(\mathbf{w}_{t_{K+1}}) - f(\mathbf{w}_{t_0}) > -\infty$, $\lim_{k \rightarrow \infty} \mathbf{d}_{t_k}^T B_{t_k} \mathbf{d}_{t_k} = 0$. Thus, $U^T \mathbf{d}_{t_k} \rightarrow \mathbf{0}$. That is to say, $\lim_{k \rightarrow \infty} \mathbf{d}_{t_k} \in \mathcal{T}^\perp$. If $\mathbf{d}_t \in \mathcal{T}^\perp$, the line search procedure will always pick unit step size. And in the next iteration, $\mathbf{d}_{t+1} = \mathbf{0}$. So when $U^T \mathbf{d}_{t_k} \rightarrow \mathbf{0}$, we also have $\mathbf{d}_{t_k} \rightarrow \mathbf{0}$. Therefore, \mathbf{w}_{t_k} converges to the optimum according to Proposition 2.5 in [26].

□

B Algorithm Details

Algorithm 2 Proximal Quasi-Newton Algorithm

Input: Observations $\{\mathbf{x}^{(i)}\}_{i=1,2,\dots,N}$, labels $\{\mathbf{y}^{(i)}\}_{i=1,2,\dots,N}$, termination criterion ϵ , scalar λ and L-BFGS memory size m .

Output: \mathbf{w}^* converging to $\arg \min_{\mathbf{w}} f(\mathbf{w})$

```

1: Initialize  $\gamma = 1$ ,  $\mathbf{w} \leftarrow \mathbf{0}$ ,  $\mathbf{g} \leftarrow \partial \ell(\mathbf{w}) / \partial \mathbf{w}$ , working set  $\mathcal{A} \leftarrow \{1, 2, \dots, d\}$ ,  $\hat{M} \leftarrow \infty$ , and  $S, Y, Q, \hat{Q} \leftarrow \phi$ .
2: for  $n = 0, 1, \dots$  do
3:    $\hat{\mathcal{A}} \leftarrow \mathcal{A}$ ,  $\mathcal{A} \leftarrow \phi$ ,  $M \leftarrow 0$ 
4:   for  $j$  in  $\hat{\mathcal{A}}$  do ▷ Shrink the working set
5:     calculate  $\partial_j f$ 


$$\partial_j f(\mathbf{w}) = \begin{cases} g_j + \text{sgn}(w_j)\lambda & \text{if } w_j \neq 0 \\ \text{sgn}(g_j) \max\{|g_j| - \lambda, 0\} & \text{if } w_j = 0 \end{cases} \quad (16)$$


6:     if  $w_j \neq 0$  or  $|g_j| - \lambda + \hat{M}/N > 0$  then
7:        $\mathcal{A} \leftarrow \mathcal{A} \cup j$ ,  $M \leftarrow \max\{M, |\partial_j f|\}$ 
8:     end if
9:   end for
10:   $\hat{M} \leftarrow M$ 
11:  if Shrinking stopping criterion attained then ▷ Check shrinking stopping criterion
12:    if Stopping criterion attained and  $|\hat{\mathcal{A}}| = d$  then ▷ Check global stopping criterion
13:      return  $\mathbf{w}$ 
14:    else
15:       $\mathbf{g} \leftarrow \partial \ell(\mathbf{w}) / \partial \mathbf{w}$ ,  $\mathcal{A} \leftarrow \{1, 2, \dots, d\}$  and  $S, Y, Q, \hat{Q} \leftarrow \phi$ 
16:      Update shrinking stopping criterion and then continue
17:    end if
18:  end if
19:   $\mathbf{d} \leftarrow \mathbf{0}$ ,  $\hat{\mathbf{d}} \leftarrow \mathbf{0}$ 
20:  Compute  $inner\_iter = \min\{max\_inner, \lfloor \frac{d}{|\mathcal{A}|} \rfloor\}$ 
21:  for  $p = 1, 2, \dots, inner\_iter$  do ▷ Solve inner problem
22:    for  $j$  in  $\mathcal{A}$  do
23:       $B_{jj} = \gamma - \mathbf{q}_j^T \hat{\mathbf{q}}_j$ ,  $(B\mathbf{d})_j = \gamma d_j - \mathbf{q}_j^T \hat{\mathbf{d}}$ 
24:       $\mathbf{a} = (B_t)_{jj}$ ,  $\mathbf{b} = (\mathbf{g}_t)_j + (B_t \mathbf{d})_j$  and  $\mathbf{c} = (\mathbf{w}_t)_j + d_j$ 
25:      Compute  $z$  according to  $z = -\mathbf{c} + \mathcal{S}(\mathbf{c} - \mathbf{b}/\mathbf{a}, \lambda/\mathbf{a})$ 
26:       $d_j \leftarrow d_j + z$ ,  $\hat{\mathbf{d}} \leftarrow \hat{\mathbf{d}} + z \hat{\mathbf{q}}_j$ 
27:    end for
28:  end for
29:  for  $\alpha = \beta^0, \beta^1, \dots$  do ▷ Conduct line search
30:    if  $f(\mathbf{w} + \alpha \mathbf{d}) \leq f(\mathbf{w}) + \alpha \sigma(\lambda \|\mathbf{w} + \mathbf{d}\|_1 - \lambda \|\mathbf{w}\|_1 + \mathbf{g}^T \mathbf{d})$  then
31:      break
32:    end if
33:  end for
34:  for  $j$  in  $\mathcal{A}$  do
35:     $g_j^{new} = \partial \ell(\mathbf{w}) / \partial w_j$ ,  $y_j = g_j^{new} - g_j$ ,  $s_j = \alpha d_j$ ,  $g_j = g_j^{new}$ 
36:  end for
37:  Update  $S, Y$  and  $Q$  just on the rows corresponding to  $\mathcal{A}$ .
38:  Update  $\gamma, D, L, S^T S$  where the inner product between  $\mathbf{s}$  and another vector is computed just over  $\mathcal{A}$ .
39:  Update  $R$  and then update  $\hat{Q}$  just on the columns corresponding to  $\mathcal{A}$ .
40: end for

```

C Proof of Theorem 3

Proof. The Hessian of $\ell(\mathbf{w})$ for CRF MLEs is

$$H = \sum_{i=1}^N \left(E \left[\phi(\mathbf{y}, \mathbf{x}^{(i)}) \phi(\mathbf{y}, \mathbf{x}^{(i)})^T \right] - E \left[\phi(\mathbf{y}, \mathbf{x}^{(i)}) \right] E \left[\phi(\mathbf{y}, \mathbf{x}^{(i)}) \right]^T \right), \quad (17)$$

where $\phi(\mathbf{y}, \mathbf{x}^{(i)}) = [f_1(\mathbf{y}, \mathbf{x}^{(i)}), f_2(\mathbf{y}, \mathbf{x}^{(i)}), \dots, f_d(\mathbf{y}, \mathbf{x}^{(i)})]^T$ and E is the expectation over the conditional probability $P_{\mathbf{w}}(\mathbf{y}|\mathbf{x}^{(i)})$. Now we re-formulate (17) to

$$H = \Phi D \Phi$$

Here $D \in \mathbb{R}^{(N|\mathcal{Y}|) \times (N|\mathcal{Y}|)}$ is a diagonal matrix with diagonal elements $D_{nn} = P_{\mathbf{w}}(\mathbf{y}_l|\mathbf{x}^{(i)})$, where $n = (i-1)|\mathcal{Y}| + l$ and $l = 1, 2, \dots, |\mathcal{Y}|$. Φ is a $d \times (N|\mathcal{Y}|)$ matrix whose column n is defined as $\Phi_n = \phi(\mathbf{y}_l, \mathbf{x}^{(i)}) - E[\phi(\mathbf{y}, \mathbf{x}^{(i)})]$ for $n = (i-1)|\mathcal{Y}| + l$.

The theorem holds because of the following four reasons.

a. \mathcal{N} is constant with respect to \mathbf{w} .

\mathcal{N} is equivalent to

$$\mathcal{N} = \{\mathbf{a} \in \mathbb{R}^d | \forall i, \exists \text{ some constant } b_i, \langle \mathbf{a}, \phi(\mathbf{y}, \mathbf{x}^{(i)}) \rangle = b_i \text{ for } \forall \mathbf{y}\} \quad (18)$$

Thus \mathcal{N} is independent on \mathbf{w} and so is \mathcal{T} .

b. $\ell(\mathbf{w})$ depends only on $\mathbf{z} = \text{proj}_{\mathcal{T}}(\mathbf{w})$.

Let $\mathbf{w} = \mathbf{z} + \mathbf{u}$. So $\mathbf{u} \in \mathcal{N}$.

$$\begin{aligned} P_{\mathbf{w}}(\mathbf{y}^{(i)}|\mathbf{x}^{(i)}) &= \frac{\exp \{ \langle \mathbf{w}, \phi(\mathbf{y}^{(i)}, \mathbf{x}^{(i)}) \rangle \}}{\sum_{\mathbf{y}} \exp \{ \langle \mathbf{w}, \phi(\mathbf{y}, \mathbf{x}^{(i)}) \rangle \}} \\ &= \frac{\exp \{ \langle \mathbf{z}, \phi(\mathbf{y}^{(i)}, \mathbf{x}^{(i)}) \rangle \} \exp \{ \langle \mathbf{u}, \phi(\mathbf{y}^{(i)}, \mathbf{x}^{(i)}) \rangle \}}{\sum_{\mathbf{y}} \exp \{ \langle \mathbf{z}, \phi(\mathbf{y}, \mathbf{x}^{(i)}) \rangle \} \exp \{ \langle \mathbf{u}, \phi(\mathbf{y}, \mathbf{x}^{(i)}) \rangle \}} \\ &= \frac{\exp \{ \langle \mathbf{z}, \phi(\mathbf{y}^{(i)}, \mathbf{x}^{(i)}) \rangle \}}{\sum_{\mathbf{y}} \exp \{ \langle \mathbf{z}, \phi(\mathbf{y}, \mathbf{x}^{(i)}) \rangle \}} \end{aligned}$$

The last equality comes from the character of \mathcal{N} , Equation (18).

c. The first property Eq. (4) holds.

$D_{nn} \rightarrow 0$ iff $\|\mathbf{w}\|_1 \rightarrow \infty$ which is prohibited by ℓ_1 penalty. Thus there exists $m_p > 0$ such that $D_{nn} \geq m_p$ for any n . Hence, the positive definiteness of H is determined by Φ .

So we have for any $\mathbf{v} \in \mathcal{T}$,

$$m_p \lambda_{\min}(\Phi \Phi^T) \|\mathbf{v}\|^2 \leq m_p \mathbf{v}^T \Phi \Phi^T \mathbf{v} \leq \mathbf{v}^T H \mathbf{v} \leq \mathbf{v}^T \Phi \Phi^T \mathbf{v} \leq \lambda_{\max}(\Phi \Phi^T) \|\mathbf{v}\|^2$$

where $\lambda_{\min}(\Phi \Phi^T)$ is the minimum nonzero eigenvalue of $\Phi \Phi^T$ and $\lambda_{\max}(\Phi \Phi^T)$ is the maximum eigenvalue of $\Phi \Phi^T$.

d. The second property Eq. (5) holds.

This property directly follows the definition of \mathcal{N} . □

D Gradient evaluation in sequence labeling and hierarchical classification

The gradients for general CRF problems are given by

$$\frac{\partial \ell(\mathbf{w})}{\partial w_k} = \sum_{i=1}^N \left(\sum_{\mathbf{y} \in \mathcal{Y}} P_{\mathbf{w}}(\mathbf{y}|\mathbf{x}^{(i)}) f_k(\mathbf{y}, \mathbf{x}^{(i)}) - f_k(\mathbf{y}^{(i)}, \mathbf{x}^{(i)}) \right) \quad (19)$$

D.1 Sequence labeling

The partial gradients of $\ell(\mathbf{w})$ for sequence labeling problem are,

$$\frac{\partial \ell(\Theta, \Lambda)}{\partial \Theta_{y,j}} = \sum_{i=1}^N \sum_{t=1}^{T^{(i)}} \left(P_{\mathbf{w}}(y_t = y | \mathbf{x}^{(i)}) - \mathbf{1} [y_t^{(i)} = y] \right) x_{tj}^{(i)} \quad (20)$$

$$\frac{\partial \ell(\Theta, \Lambda)}{\partial \Lambda_{y,y'}} = \sum_{i=1}^N \sum_{t=1}^{T^{(i)}-1} \left(P_{\mathbf{w}}(y_t = y, y_{t+1} = y' | \mathbf{x}^{(i)}) - \mathbf{1} [y_t^{(i)} = y, y_{t+1}^{(i)} = y'] \right) \quad (21)$$

The forward-backward algorithm is a popular inference oracle for evaluating the marginal probability in Equation (20) and (21). In our OCR model, the forward-backward algorithm is

$$\begin{cases} \alpha_1(y) = \exp(\Theta_y^T \mathbf{x}_1) \\ \alpha_{t+1}(y) = \sum_{y'} \alpha_t(y') \exp(\Theta_y^T \mathbf{x}_{t+1} + \Lambda_{y',y}) \\ \beta_T(y) = 1 \\ \beta_t(y') = \sum_y \beta_{t+1}(y) \exp(\Theta_y^T \mathbf{x}_{t+1} + \Lambda_{y',y}) \end{cases}$$

where Θ_y^T is the y -th row of the matrix Θ . Then the marginal conditional probabilities are given by

$$\begin{aligned} P_{\mathbf{w}}(y_t = y', y_{t+1} = y | \mathbf{x}) &= \frac{1}{Z_{\mathbf{w}}(\mathbf{x})} \alpha_t(y') \exp(\Theta_y^T \mathbf{x}_{t+1} + \Lambda_{y',y}) \beta_{t+1}(y) \\ P_{\mathbf{w}}(y_t = y | \mathbf{x}) &= \frac{1}{Z_{\mathbf{w}}(\mathbf{x})} \alpha_t(y) \beta_t(y), \end{aligned}$$

where the normalization factor $Z_{\mathbf{w}}(\mathbf{x})$ can be computed by $\sum_y \alpha_T(y)$.

D.2 Hierarchical classification

The partial gradients of $\ell(W)$ for hierarchical classification problem are,

$$\frac{\partial \ell(W)}{\partial W_{k,j}} = \sum_{i=1}^N \left(\sum_{y \in \mathcal{Y}} \mathbf{1} [k \in \text{Path}(y)] P_W(y | \mathbf{x}^{(i)}) - \mathbf{1} [k \in \text{Path}(y^{(i)})] \right) x_j^{(i)}$$

They can be evaluated by the downward-upward algorithm. Let $\alpha(k)$ and $\beta(k)$ be the downward message and upward message respectively.

$$\begin{cases} \alpha(\text{root}) = \mathbf{w}_{\text{root}}^T \mathbf{x} \\ \alpha(k) = \alpha(\text{parent}(k)) + \mathbf{w}_k^T \mathbf{x} \\ \beta(k) = \alpha(k) / \sum_{y \in Y} \alpha(y) & \text{if } k \text{ is a leaf node} \\ \beta(k) = \sum_{k' \in \text{children}(k)} \beta(k') & \text{if } k \text{ is a non-leaf node} \end{cases}$$

So we have

$$\frac{\partial \ell(W)}{\partial W_{k,j}} = \sum_{i=1}^N \left(\beta^{(i)}(k) - \mathbf{1} [k \in \text{Path}(y^{(i)})] \right) x_j^{(i)} \quad (22)$$

E More Experimental Results

E.1 Performance on different values of λ

λ affects the sparsity of the intermediate iterates, and further the speed of the algorithm. In particular, when λ is larger, the intermediate iterates are sparser, and then the corresponding iterations, due to the shrinking strategy, will be faster – and vice versa. The effect of λ on the performance is shown in this section.

E.1.1 Sequence Labeling

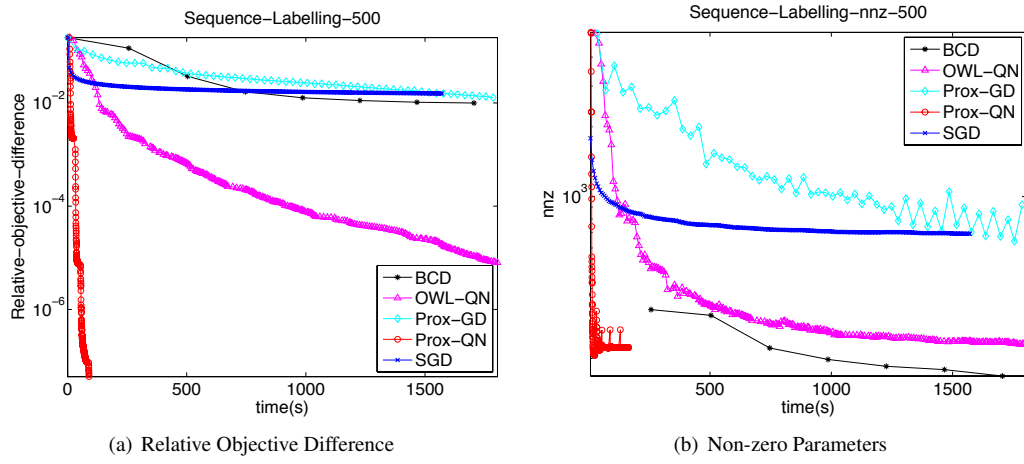


Figure 3: Sequence Labeling Problem for $\lambda = 500$

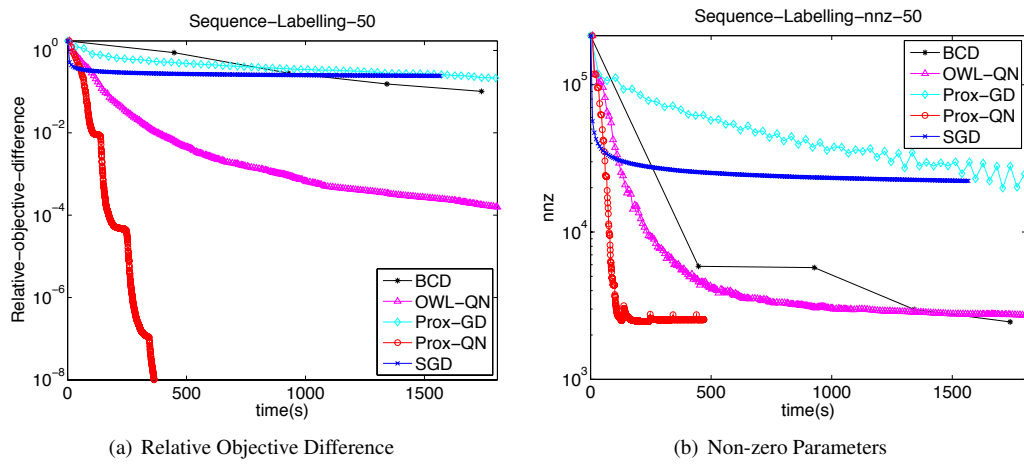


Figure 4: Sequence Labeling Problem for $\lambda = 50$

E.1.2 Hierarchical Classification

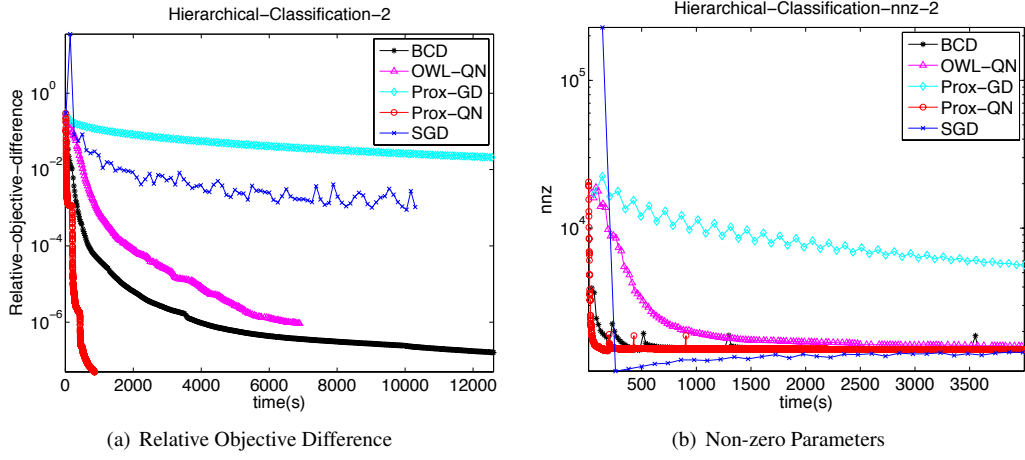


Figure 5: Hierarchical Classification Problem for $\lambda = 2$

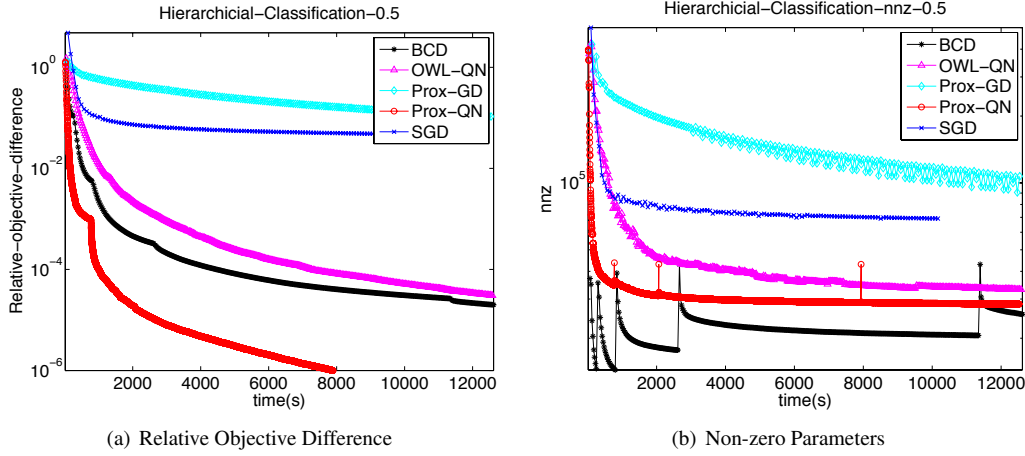


Figure 6: Hierarchical Classification Problem for $\lambda = 0.5$

E.2 Testing Accuracy

The testing accuracy for different λ 's on these two problems is in the following tables. The testing accuracy across training time is shown in Figure 7.

| λ | 50 | 100 | 500 |
|------------------|----------|----------|----------|
| Testing Accuracy | 0.834928 | 0.736643 | 0.407895 |
| nnz of optimum | 2542 | 1544 | 223 |

Table 1: Per-character testing accuracy for OCR dataset

| λ | 0.5 | 1 | 2 |
|------------------|----------|----------|----------|
| Testing Accuracy | 0.262648 | 0.249731 | 0.185684 |
| nnz of optimum | 28483 | 4301 | 1505 |

Table 2: Testing accuracy for LSHTC1 dataset

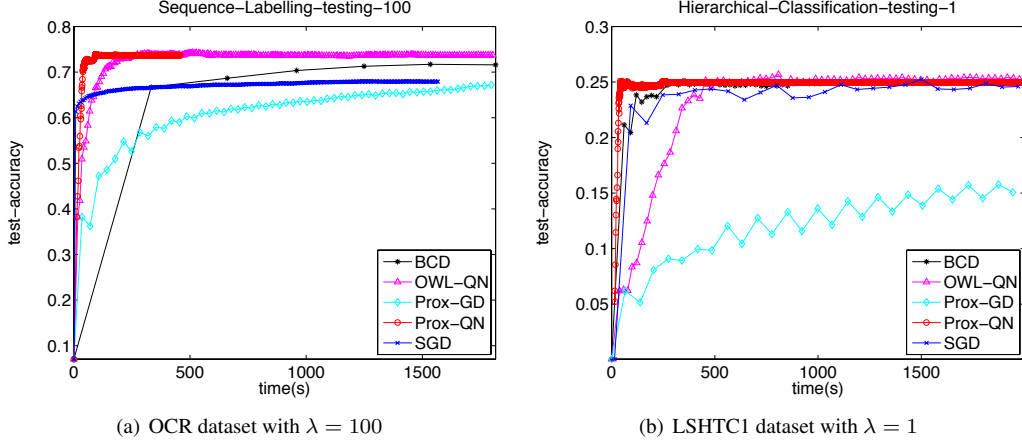


Figure 7: Testing accuracy v.s. training time

E.3 Relative objective difference v.s. the number of passes over dataset

Figure 8 shows the performance under the measure of number of passes (iterations) over dataset. We do not include the plot for BCD method, because the pass over dataset for BCD actually depends on the sparsity pattern of the dataset. Thus it is hard to fairly define the pass over dataset for BCD. In this experiment, shrinking strategy is not applied.

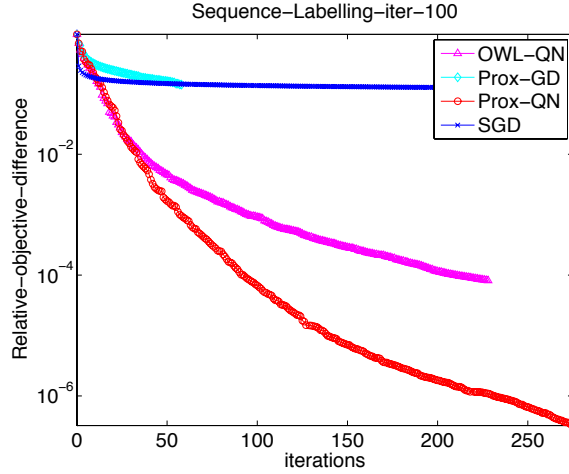


Figure 8: Relative objective difference v.s. the number of passes over dataset for OCR dataset with $\lambda = 100$.